

# Multi-Task Triplet Loss for Named Entity Recognition using Supplementary Text

Ryan Siskind {ryan.siskind@target.com}  
Shalin Shah {shalin.shah@target.com}

AI Sciences  
Target Corporation  
Sunnyvale, CA 94086, USA

## Abstract

Retail item data contains many different forms of text like the title of an item, the description of an item, item name and reviews. It is of interest to identify the item name in the other forms of text using a named entity tagger. However, the title of an item and its description are syntactically different (but semantically similar) in that the title is not necessarily a well formed sentence while the description is made up of well formed sentences. In this work, we use a triplet loss to contrast the embeddings of the item title with the description to establish a proof of concept. We find that using the triplet loss in a multi-task NER algorithm improves both the precision and recall by a small percentage. While the improvement is small, we think it is a step in the right direction of using various forms of text in a multi-task algorithm. In addition to precision and recall, the multi-task triplet loss method is also found to significantly improve the exact match accuracy i.e. the accuracy of tagging the entire set of tokens in the text with correct tags.

*Keywords:* named entity recognition, NER, multi-task learning, contrastive loss, triplet loss, BERT, neural networks, deep learning, supplementary text

## 1 Introduction

Bayesian Personalized Ranking [1] [2] was one of the initial methods that used the triplet loss for personalized recommender systems. The triplet loss [3] is similar to the contrastive loss [4] [5] [6] but uses a triplet instead of two sets of embeddings. In this work, we use the triplet loss in a multi-task learning method to learn named entity recognition. An item in the catalog of Target has several types of textual data available that in some way describe the item.

There is a title, which is a short description of the item, and there is a description which is a long textual representation of the item. Instead of doing named entity recognition by adding the title and the description as separate rows, we use them together through a triplet loss on the embeddings. We find that this improves the results as compared to BERT-base [7].

Contrastive loss and the triplet loss are similar in that they both contrast the record under consideration with other records. In a contrastive loss, a row is learned to be similar to another positive row and it is learned to be not similar to another negative row. In a triplet loss, the objective is to maximize the difference between the similarity of the row under consideration with the positive and negative rows respectively.

The goal in this paper is to use supplementary text (like item descriptions) with the more important item titles in a multi-task triplet loss.

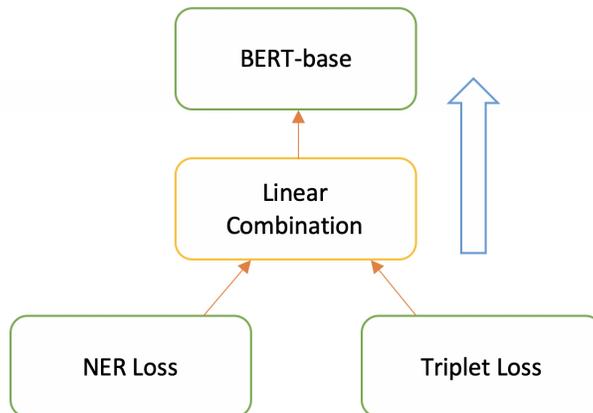


Figure 1: Triplet Loss and NER Loss using BERT-base

## 2 Methods

We use the triplet loss to contrast the similarities and differences between BERT [7] sentence embeddings of titles and descriptions. This loss is very similar to the loss used in BPR [1] except that we have two objectives in our loss function and the goal is to learn efficient sentence embeddings as opposed to user and item latent factors.

Let  $t_i$  be the title embedding of the  $i^{th}$  title and let  $d_p$  and  $d_n$  be the sentence embeddings two descriptions, where  $d_p$  is the description of the  $i^{th}$  item

under consideration and  $d_n$  is a randomly chosen description of a negative item.

Then, the triplet loss is the following:

$$c_p = \text{cosine}(t_i, d_p)$$

$$c_n = \text{cosine}(t_i, d_n)$$

$$d_i = c_p - c_n \text{ (where the objective is to maximize)}$$

$$\mathcal{L} = \frac{e^{d_i}}{1+e^{d_i}}$$

We use this loss  $\mathcal{L}$  in a multi-task setting with the named entity recognition loss. We propagate both losses backwards in the neural network. However, during scoring (inference), the descriptions are not used at all.

We use BERT-base [7] as the common algorithm and we build upon BERT-base through transfer learning. We do not lock any weights during backpropagation.

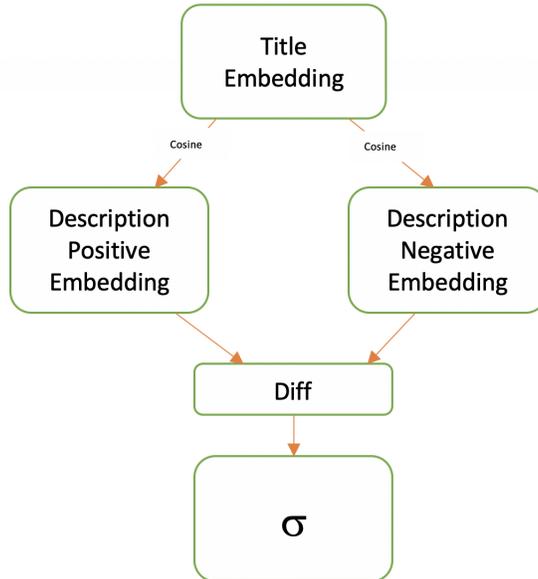


Figure 2: Title-Description Triplet Loss

Table 1: Results on a 30% held out test set (same random seeds)

Algorithm	Precision	Recall	Exact Matches	Accuracy
BERT-Multitask-Triplet	<b>78%</b>	<b>63%</b>	<b>43%</b>	<b>85%</b>
BERT-base	77%	62%	41%	84.7%

### 3 Results

Table 1 shows the results of our BERT-Multitask-Triplet. As the table shows the BERT-Multitask-Triplet algorithm improves upon precision and recall by a percentage point.

An exact match is a labeling of the sentence in which all tags are correct. As the results show, the BERT-Multitask-Triplet method is able to improve the exact match accuracy by two percentage points which is quite significant.

The improvement in overall accuracy is marginal, but still relevant to show that the BERT-Multitask-Triplet outperforms BERT-base.

### 4 Conclusion

In this work, we presented an algorithm BERT-Multitask-Triplet for named entity recognition (NER) and find that after transfer learning, the algorithm is able to outperform BERT-base. Our algorithm is a multi-objective learning algorithm which linearly combines the NER objective with the triplet loss objective. The triplet loss objective uses supplementary text data of the items in the catalog of Target like the item descriptions.

This method is applicable to many other domains like in medicine where notes of multiple doctors might be available for a patient and the goal is to some form of classification. Instead of concatenating the text, the supplementary text could be used in a multi-objective learning algorithm which contrasts the sentence embeddings using a loss as described above.

It could be useful to compare algorithms using BERT-large to see if the improvements carry forward as the number of parameters of the neural network increase. It might be useful to compare the algorithm with contrastive losses to see if the triplet loss does indeed work better.

## References

- [1] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. Bpr: Bayesian personalized ranking from implicit feedback. *arXiv preprint arXiv:1205.2618*, 2012.
- [2] Shalin Shah. A survey of latent factor models for recommender systems and personalization. 2021.
- [3] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015.
- [4] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *arXiv preprint arXiv:2004.11362*, 2020.
- [5] Christopher JC Burges. From ranknet to lambdarank to lambdamart: An overview. *Learning*, 11(23-581):81, 2010.
- [6] Beliz Gunel, Jingfei Du, Alexis Conneau, and Ves Stoyanov. Supervised contrastive learning for pre-trained language model fine-tuning. *arXiv preprint arXiv:2011.01403*, 2020.
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.